

CLAIMS

1 1. A method of extracting relevant data, comprising:
2 accessing at least a first set of data of a first tree, wherein the first set of data
3 includes selected data of the first tree, the selected data at least partly specifying tree
4 data;
5 accessing at least a second set of data of a second tree;
6 determining an edit sequence between at least part of the first set of data and at
7 least part of the second set of data, the edit sequence including any of insertions,
8 deletions, substitutions, matches, and repetitions; and
9 finding corresponding data of the second set of data, the corresponding data
10 having a correspondence to the selected data, the correspondence at least partly found
11 by determining the edit sequence.

1 2. The method of claim 1, wherein the repetitions are subtrees of at least
2 the first tree.

1 3. The method of claim 1, wherein the edit sequence includes at least two
2 repetitions, the at least two repetitions based on at least one subtree of the first tree, and
3 the at least two repetitions appears in the second tree, and the at least two repetitions
4 include at least a first repetition and a second repetition, and the first repetition has at
5 least one difference from the second repetition.

1 4. The method of claim 3, wherein each of the at least two repetitions is
2 obtainable from the at least one subtree of the first tree by some sequence of one or
3 more insertions, deletions, substitutions and matches.

1 5. The method of claim 1, wherein the edit sequence includes none of
2 insertions, deletions, substitutions, matches, and repetitions.

1 6. The method of claim 1, wherein the edit sequence includes at least one
2 of one or more insertions of nodes, one or more insertions of subtrees, one or more
3 deletions of subtrees, one or more deletions of subtrees, one or more substitutions of
4 nodes, one or more substitutions of subtrees, one or more repetitions of nodes, and one
5 or more repetitions of subtrees.

1 7. The method of claim 1, wherein the edit sequence is at least partly
2 determined by calculating a total cost, and each of one or more of insertions, deletions,
3 substitutions, and matches is associated with one or more costs.

1 8. The method of claim 7, wherein the one or more costs are at least
2 partly set to encourage the edit sequence to include one or more matches between at
3 least some selected data of the first tree and at least some data from the second tree.

1 9. The method of claim 7, wherein the one or more costs are at least
2 partly set to encourage the edit sequence to include one or more repetitions.

1 10. The method of claim 7, wherein a first cost is associated with a first
2 match at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second match at a second distance from a root of a tree
4 representation of some set of data, the first distance is less than the second distance,
5 and the first cost and the second cost are different.

1 11. The method of claim 7, wherein a first cost is associated with a first
2 insertion at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second insertion at a second distance from a root of a
4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 12. The method of claim 7, wherein a first cost is associated with a first
2 deletion at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second deletion at a second distance from a root of a
4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 13. The method of claim 7, wherein a first cost is associated with a first
2 substitution at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second substitution at a second distance from a root of
4 a tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 14. The method of claim 7, wherein a first cost is associated with a first
2 repetition at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second repetition at a second distance from a root of a
4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 15. The method of claim 7, wherein a first cost is associated with a first
2 text-based content substitution such that a first length of substituting text-based content
3 is substantially equal to a first length of substituted text-based content, a second cost is
4 associated with a second text-based content substitution such that a second length of
5 substituting text-based content is substantially different from a second length of
6 substituted text-based content, and the first cost and the second cost are set to
7 discourage the second text-based content substitution more than the first text-based
8 content substitution.

1 16. The method of claim 7, wherein data includes at least a first type and a
2 second type, and the one or more costs are at least partly set to discourage substitutions
3 of one or more of the first type for one or more of the second type.

1 17. The method of claim 7, wherein data includes at least a first type and a
2 second type, and the one or more costs are at least partly set to discourage substitutions
3 of one or more of the second type for one or more of the first type.

1 18. The method of claim 7, wherein a first cost is associated with
2 preserving data of a first type with unchanged attributes, a second cost is associated
3 with preserving data of a second type with one or more changed attributes, and the first
4 cost and the second cost are set to discourage preserving data of the second type more
5 than preserving the data of the first type.

1 19. The method of claim 1, wherein tree data is at least partly from the first
2 tree.

1 20. The method of claim 1, wherein tree data is at least partly from the
2 second tree.

1 21. The method of claim 1, wherein the second tree is received if the
2 second tree is different from the first tree.

1 22. The method of claim 1, further comprising:
2 if two or more corresponding data are found, then:
3 selecting larger selected data, at least part of the larger selected data
4 including a larger subtree in a first tree representation of the first set of data, the larger
5 subtree including the selected data;
6 determining a second edit sequence between at least part of the first set
7 of data and at least part of a second tree representation of the second set of data, the
8 first set of data including at least part of the larger selected data, the second edit
9 sequence including any of insertions, deletions, and substitutions;
10 finding corresponding data of the second set of data, the corresponding
11 data having a correspondence to the larger selected data, the correspondence at least
12 partly found by determining the second edit sequence; and
13 finding corresponding data of the second set of data, the corresponding
14 data having a correspondence to the selected data, the correspondence at least partly
15 found by determining the second edit sequence.

1 23. The method of claim 1, wherein the correspondence is at least partly
2 found by one or more of: determining the edit sequence, at least part of at least one of a
3 first plurality of paths from a root of a tree representation of the first set of data to
4 selected data of the tree representation of the first set of data, at least part of at least one
5 of a second plurality of paths from a root of a tree representation of the second set of
6 data to corresponding data of the tree representation of the second set of data, and one
7 or more edit sequences between at least one of the first plurality of paths and at least
8 one of the second plurality of paths.

1 24. The method of claim 1, wherein one or more of the first set of data and
2 the second set of data is represented at least partly by a tree.

1 25. The method of claim 1, wherein one or more of the first set of data and
2 the second set of data is represented at least partly by a set of linearized tokens.

1 26. The method of claim 1, wherein the first tree and the second tree
2 represent different trees.

1 27. The method of claim 1, wherein the first tree and the second tree
2 represent a same tree.

1 28. The method of claim 1, wherein the first tree and the second tree
2 represent different versions of a same tree.

1 29. The method of claim 1, further comprising:
2 determining at least one edit sequence of forward and backward edit
3 sequences between at least part of the first tree and at least part of the second tree;

1 performing at least one of 1) and 2):

2 1a) pruning a relevant subtree from at least part of the first tree,
3 the relevant subtree at least partly determined from the forward and backward edit
4 sequences;

5 1b) determining a pruned edit sequence between the pruned
6 relevant subtree and at least part of the second tree;

7 2a) pruning a relevant subtree from at least part of the second tree,
8 the relevant subtree at least partly determined from the forward and backward edit
9 sequences;

10 2b) determining a pruned edit sequence between at least part of the
11 first tree and the pruned relevant subtree; and

12 finding corresponding data of the second set of data, the corresponding
13 data having a correspondence to the selected data, the correspondence at least partly
14 found by determining the pruned edit sequence.

1 30. A method of extracting relevant data, comprising:
2 accessing at least a first set of data of a first tree, wherein the first set
3 of data includes selected data of the first tree, the selected data at least partly specifying
4 tree data;

5 accessing at least a second set of data of a second tree;

6 determining a second path from a root of the second tree that
7 corresponds to a first path from a root of the first tree to the selected data; and

8 finding corresponding data of the second set of data, the corresponding
9 data having a correspondence to the selected data, the correspondence at least partly
10 determined by the second path from the root of the second tree.

1 31. The method of claim 30, wherein the second path is determined at least
2 in part by:
3 traversing the first tree and the second tree;
4 at each traversed level of the first tree, the traversed level of the first
5 tree including a plurality of level nodes, selecting a level node of the plurality of level
6 nodes, the level node being in the first path;
7 at each traversed level of the second tree, selecting a best
8 corresponding node at the traversed level of the second tree, the best corresponding
9 node saving a best correspondence to the selected level node of the plurality of level
10 nodes.

1 32. The method of claim 31, wherein the best corresponding node is
2 determined at least in part by determining an edit sequence between a first subset of
3 data obtained from at least part of the first set of data and a second subset of data
4 obtained from at least part of the second set of data, the edit sequence including any of
5 insertions, deletions, substitutions, matches, and repetitions.

1 33. The method of claim 32, wherein the repetitions are subtrees of at least
2 the first tree.

1 34. The method of claim 32, wherein the edit sequence includes at least
2 two repetitions, the at least two repetitions based on at least one subtree of the first tree,
3 and the at least two repetitions appears in the second tree, and the at least two
4 repetitions include at least a first subtree and a second subtree, and the first subtree has
5 at least one difference from the second subtree.

1 35. The method of claim 34, wherein each of the at least two repetitions is
2 obtainable from the at least one subtree of the first tree by some sequence of one or
3 more insertions, deletions, substitutions and matches.

1 36. The method of claim 32, wherein the edit sequence includes none of
2 insertions, deletions, substitutions, matches, and repetitions.

1 37. The method of claim 32, wherein the edit sequence includes at least
2 one of one or more insertions of nodes, one or more insertions of subtrees, one or more
3 deletions of subtrees, one or more deletions of subtrees, one or more substitutions of
4 nodes, one or more substitutions of subtrees, one or more repetitions of nodes, and one
5 or more repetitions of subtrees.

1 38. The method of claim 32, wherein the edit sequence is at least partly
2 determined by calculating a total cost, and each of one or more of insertions, deletions,
3 substitutions, and matches is associated with one or more costs.

1 39. The method of claim 38, wherein the one or more costs are at least
2 partly set to encourage the edit sequence to include one or more matches between at
3 least some selected data of the first tree and at least some data from the second tree.

1 40. The method of claim 38, wherein the one or more costs are at least
2 partly set to encourage the edit sequence to include one or more repetitions.

1 41. The method of claim 38, wherein a first cost is associated with a first
2 match at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second match at a second distance from a root of a tree
4 representation of some set of data, the first distance is less than the second distance,
5 and the first cost and the second cost are different.

1 42. The method of claim 38, wherein a first cost is associated with a first
2 insertion at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second insertion at a second distance from a root of a
4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 43. The method of claim 38, wherein a first cost is associated with a first
2 deletion at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second deletion at a second distance from a root of a

4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 44. The method of claim 38, wherein a first cost is associated with a first
2 substitution at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second substitution at a second distance from a root of
4 a tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 45. The method of claim 38, wherein a first cost is associated with a first
2 repetition at a first distance from a root of a tree representation of some set of data, a
3 second cost is associated with a second repetition at a second distance from a root of a
4 tree representation of some set of data, the first distance is less than the second
5 distance, and the first cost and the second cost are different.

1 46. The method of claim 38, wherein a first cost is associated with a first
2 text-based content substitution such that a first length of substituting text-based content
3 is substantially equal to a first length of substituted text-based content, a second cost is
4 associated with a second text-based content substitution such that a second length of
5 substituting text-based content is substantially different from a second length of
6 substituted text-based content, and the first cost and the second cost are set to
7 discourage the second text-based content substitution more than the first text-based
8 content substitution.

1 47. The method of claim 38, wherein data includes at least a first type and
2 a second type, and the one or more costs are at least partly set to discourage
3 substitutions of one or more of the first type for one or more of the second type.

1 48. The method of claim 38, wherein data includes at least a first type and
2 a second type, and the one or more costs are at least partly set to discourage
3 substitutions of one or more of the second type for one or more of the first type.

1 49. The method of claim 38, wherein a first cost is associated with
2 preserving data of a first type with unchanged attributes, a second cost is associated
3 with preserving data of a second type with one or more changed attributes, and the first

4 cost and the second cost are set to discourage preserving data of the second type more
5 than preserving the data of the first type.

1 50. The method of claim 32, wherein the first subset of data includes nodes
2 in the first tree that are within a first neighborhood of any of the selected level nodes of
3 the traversed levels of the first tree, the selected level nodes of the traversed levels of
4 the first tree being on the first path from the root of the first tree to the selected data,
5 and the second subset of data includes nodes in the second tree that are within a second
6 neighborhood of children nodes of the best corresponding node selected at a previous
7 level in the second tree.

1 51. The method of claim 50, wherein the first neighborhood of any
2 selected level node includes a first plurality of close nodes according to a first distance
3 measure, and the second neighborhood of a child node includes a second plurality of
4 close nodes according to a second distance measure.

1 52. The method of 51, wherein the first distance measure between any
2 selected level node and another node is at least partly determined a first number of tree
3 edges between any selected level node and another node the second distance measure
4 between the child node and another node is at least partly determined a second number
5 of tree edges between the child node and another node.

1 53. The method of 51, wherein the first distance measure between any
2 selected level node and another node is at least partly determined a first number of tree
3 levels between any selected level node and another node the second distance measure
4 between the child node and another node is at least partly determined a second number
5 of tree levels between the child node and another node.

1 54. The method of claim 30, wherein one or more of the first set of data
2 and the second set of data is represented at least partly by a tree.

1 55. The method of claim 30, wherein one or more of the first set of data
2 and the second set of data is represented at least partly by a set of linearized tokens.

1 56. The method of claim 30, wherein the first tree and the second tree
2 represent different trees.

1 57. The method of claim 30, wherein the first tree and the second tree
2 represent a same tree.

1 58. The method of claim 30, wherein the first tree and the second tree
2 represent different versions of a same tree.

1 59. The method of claim 30, wherein tree data is at least partly from the
2 first tree.

1 60. The method of claim 30, wherein tree data is at least partly from the
2 second tree.

1 61. The method of claim 30, wherein the second tree is received if the
2 second tree is different from the first tree.

1 62. An apparatus for extracting relevant data, comprising:
2 a plurality of one or more computing devices adapted to perform:
3 accessing at least a first set of data of a first tree, wherein the first set
4 of data includes selected data of the first tree, the selected data at least partly specifying
5 tree data;
6 accessing at least a second set of data of a second tree;
7 determining an edit sequence between at least part of the first set of
8 data and at least part of the second set of data, the edit sequence including any of
9 insertions, deletions, substitutions, matches, and repetitions; and
10 finding corresponding data of the second set of data, the corresponding
11 data having a correspondence to the selected data, the correspondence at least partly
12 found by determining the edit sequence.

1 63. An apparatus for extracting relevant data, comprising:
2 a plurality of one or more computing devices adapted to perform:
3 accessing at least a first set of data of a first tree, wherein the first set
4 of data includes selected data of the first tree, the selected data at least partly specifying

5 tree data;
6 accessing at least a second set of data of a second tree;
7 determining a second path from a root of the second tree that
8 corresponds to a first path from a root of the first tree to the selected data; and
9 finding corresponding data of the second set of data, the corresponding
10 data having a correspondence to the selected data, the correspondence at least partly
11 determined by the second path from the root of the second tree.